

AIC入門

3章 推定

4章 AIC

B1

モデルの評価, 推定

偶然を伴う現象, 複雑な現象 — ある確率分布に従う確率変数の実現値とみなす

(1) データを生成する機構 (2) 現象を表現するモデル
ともに確率分布であるとみなす立場をとる

- 真の確率分布を近似する「モデル」の近似の良さの評価の方法
モデルが確率分布に対してどの程度近いかを評価する
- 与えられたデータからモデルを推定する方法
なるべく近似の良いモデル

B3

3 推定

この章のトピック — 最尤法

- エントロピーと情報量
- 情報量の推定値 — 対数尤度
- 最尤法

B2

3.1 エントロピーと情報量

[注意] 教科書の表記とは p と q を入れ替えてある.

m 個の事象 $\omega_1, \dots, \omega_m$ の確率を q_1, \dots, q_m とする離散分布を考える

$$\mathbf{q} = (q_1, \dots, q_m); \quad q_1 + \dots + q_m = 1, q_i > 0$$

[例] ある球団の勝つ確率

予想 A 7割: $p_A = (0.7, 0.3)$

予想 B 5割: $p_B = (0.5, 0.5)$

もし真の勝率が4割なら $\mathbf{q} = (0.4, 0.6)$ なので, 予想 B のほうが予想 A より良い.

もし真の勝率が6割なら $\mathbf{q} = (0.6, 0.4)$ だが, どちらがよいか?

真の分布とモデルの「近さ」を測る客観的な規準が必要

B4

K-L 情報量 (相対エントロピー)

真の離散分布 $q = (q_1, \dots, q_m)$, モデル $p = (p_1, \dots, p_m)$

$$I(q; p) = \sum_{i=1}^m q_i \log \frac{q_i}{p_i} \quad \log \text{は } e \text{ を底とする自然対数.}$$

モデル p に関する真の分布 q の Kullback-Leibler (カルバック・ライブラ) 情報量.

K-L 情報量を分布の近さを測る規準として採用する

$I(q; p)$ の値が小さく 0 に近いほどモデル p は真の分布 q に近いとみなす

[例] 勝率予想

$$I(q; p_A) = 0.6 \log \frac{0.6}{0.7} + 0.4 \log \frac{0.4}{0.3} = 0.0226$$

$$I(q; p_B) = 0.6 \log \frac{0.6}{0.5} + 0.4 \log \frac{0.4}{0.5} = 0.0201$$

予想 **B** の方が (わずかに) 良い

B5

K-L 情報量の解釈

$n = (n_1, \dots, n_m)$ が多項分布に従うとする. $\sum_i n_i = n$ で確率ベクトルは p とする.

nq の要素が整数のとき,

$$\log P\left(\frac{n}{n} = q \mid p\right) \approx -nI(q; p)$$

[参考] $-I(q; p)$ は負のエントロピーと呼ばれる.

B7

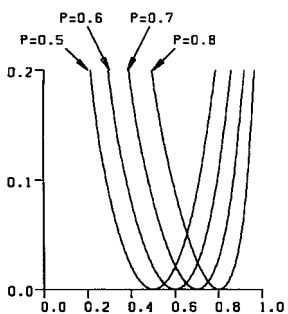
K-L 情報量の性質

(i) $I(q; p) \geq 0$

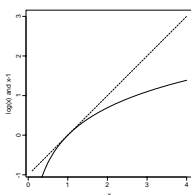
(ii) $I(q; p) = 0 \Leftrightarrow p = q \quad (p_i = q_i, i = 1, \dots, m)$

[証明] $\log x \leq x - 1$ で等号は $x = 1$ のときのみ. $x = p_i/q_i$ を代入して

$$-I(q; p) = \sum_{i=1}^m q_i \log \frac{p_i}{q_i} \leq \sum_{i=1}^m q_i \left(\frac{p_i}{q_i} - 1 \right) = \sum_{i=1}^m (p_i - q_i) = 0$$



$\log(x)$ と $x - 1$ のプロット



$q = 0.5, 0.6, 0.7, 0.8$ と変えたときの $I((q, 1 - q); (p, 1 - p))$ の p に対するプロット.

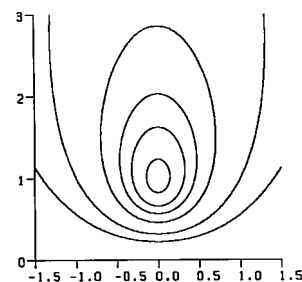
B6

真の密度関数 $g(x)$, モデル $f(x)$

$$I(g; f) = \int_{-\infty}^{\infty} g(x) \log \frac{g(x)}{f(x)} dx$$

[例] 真の分布 $g: N(\mu, \sigma^2)$, モデル $f: N(\xi, \tau^2)$

$$g(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad f(x) = \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{(x-\xi)^2}{2\tau^2}\right)$$



$$I(g; f) = \frac{1}{2} \left\{ \log \frac{\tau^2}{\sigma^2} + \frac{\sigma^2 + (\mu - \xi)^2}{\tau^2} - 1 \right\}$$

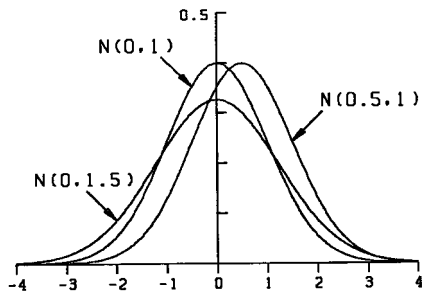
$(\mu, \sigma^2) = (0, 1)$ のときの $I(g; f)$ の (ξ, τ^2) に対するプロット

B8

[例] つぎの二つのモデルのうちどちらが真の分布 $N(0, 1)$ に近いか？

$f_1 : N(0.5, 1)$ の K-L 情報量は $I(g; f_1) = 0.125$

$f_2 : N(0, 1.5)$ の K-L 情報量は $I(g; f_2) = 0.036$



B9

3.2 情報量の推定値 — 対数尤度

実際の応用では真の分布 q は未知．(連続型の場合は $g(\cdot)$ が未知)．

データから K-L 情報量を推定する

[例] ある球団の勝つ確率

予想 A 7割: $p_A = (0.7, 0.3)$

予想 B 5割: $p_B = (0.5, 0.5)$

100 試合終わった段階での実際の成績は 65 勝 35 敗だった．

どちらの予想が良かったか？

B11

[参考] 情報理論におけるシャノンのエントロピー

$$H(q) = \sum_{i=1}^m q_i \log \frac{1}{q_i}$$

対数の底が e ならば単位は nat (ナット), 底が 2 ならば単位は bit (ビット) となる．

シャノンのエントロピーと K-L 情報量の関係

$$H(q) = H(p_0) - I(q; p_0)$$

ただし $p_0 = (1/m, \dots, 1/m)$, $H(p_0) = \log m$.

[参考] 情報理論における相互情報量

$$H(X; Y) = \sum_{x_i} \sum_{y_j} p(x_i, y_j) \log \frac{p(x_i, y_j)}{p(x_i)p(y_j)}$$

相互情報量と K-L 情報量の関係

$$H(X; Y) = I(p(x, y); p(x)p(y))$$

B10

平均対数尤度

K-L 情報量の式を展開すると

$$I(q; p) = \sum_{i=1}^m q_i \log \frac{q_i}{p_i} = \sum_{i=1}^m q_i \log q_i - \sum_{i=1}^m q_i \log p_i$$

この第 2 項を大きくするモデル p が良い

平均対数尤度 (の n 倍)

$$l^*(p) = n \sum_{i=1}^m q_i \log p_i$$

これを大きくする p が良いモデル

B12

対数尤度

事象 w_i を観測した回数を n_i とすると、相対頻度 n_i/n は q_i の推定になっている。

$l^*(p)$ の q_i に n_i/n を代入して次の推定値を得る。

対数尤度

$$l(p) = \sum_{i=1}^m n_i \log p_i$$

これを大きくする p は良いモデルの推定になる。

$i = 1, \dots, n$ に対して i 番目の観測値を x_i とすると

$$l(p) = \sum_{i=1}^n \log p_{x_i}$$

とも書ける。

B13

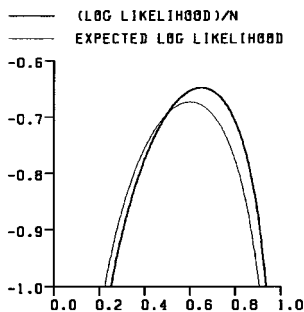
[例] 球団の勝率予想 $n_1 = 65, n_2 = 35$

$$l(p) = 65 \log p_1 + 35 \log p_2$$

$$l(p_A) = 65 \log 0.7 + 35 \log 0.3 = -65.32$$

$$l(p_B) = 65 \log 0.5 + 35 \log 0.5 = -69.31$$

データからは予想 **A** の方が良かったと判断される。



$n \rightarrow \infty$ の極限では、

$$\frac{1}{n} l(p) \rightarrow \frac{1}{n} l^*(p)$$

であるが有限の n では多少ずれる。

B14

連続型確率変数の場合

K-L 情報量

$$\begin{aligned} I(g; f) &= \int_{-\infty}^{\infty} g(x) \log \frac{g(x)}{f(x)} dx \\ &= \int_{-\infty}^{\infty} g(x) \log g(x) dx - \int_{-\infty}^{\infty} g(x) \log f(x) dx \end{aligned}$$

平均対数尤度 (の n 倍)

$$l^*(f) = n \int_{-\infty}^{\infty} g(x) \log f(x) dx$$

対数尤度

$$l(f) = \sum_{i=1}^n \log f(x_i)$$

$n \rightarrow \infty$ の極限では

$$\frac{1}{n} l(f) \rightarrow \frac{1}{n} l^*(f)$$

B15

[例] 10個の観測値 $\{-1.10, -0.40, -0.20, -0.02, 0.02, 0.71, 1.35, 1.46, 1.74, 3.89\}$ が得られたとき、つぎの二つのモデルのうちどちらが良いか？

$$f_1(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \quad \text{正規分布}$$

$$f_2(x) = \frac{1}{\pi(x^2 + 1)} \quad \text{コーシー分布}$$

$$\sum_{i=1}^{10} \log f_1(x_i) = -21.20, \quad \sum_{i=1}^{10} \log f_2(x_i) = -19.19$$

コーシー分布の方が良い。

もし最後の観測値 3.89 (外れ値?) をとりのぞくと、

$$\sum_{i=1}^9 \log f_1(x_i) = -12.72, \quad \sum_{i=1}^9 \log f_2(x_i) = -15.26$$

正規分布の方が良い。

B16

尤度

確率変数 (X_1, \dots, X_n) のパラメトリック・モデル

$$f(x_1, \dots, x_n | \theta)$$

データ (x_1, \dots, x_n) を観測したときの尤度 (**likelihood**)

$$L(\theta) = f(x_1, \dots, x_n | \theta)$$

対数尤度は、尤度の対数として定義される

$$l(\theta) = \log L(\theta)$$

X_1, \dots, X_n が独立なら、

$$L(\theta) = f(x_1 | \theta) \cdots f(x_n | \theta) \quad \text{or} \quad p(x_1 | \theta) \cdots p(x_n | \theta)$$

$$l(\theta) = \sum_{i=1}^n \log f(x_i | \theta) \quad \text{or} \quad \sum_{i=1}^n \log p(x_i | \theta)$$

B17

3.3 最尤法

パラメトリック・モデルにおいて、データが与えられたとき、尤度 $L(\theta)$ もしくは対数尤度 $l(\theta)$ を最大とするようなパラメタの値 $\hat{\theta}$ を選ぶ

このパラメタ推定法は、最大尤度法、略して「最尤法」と呼ばれる。

推定されたパラメタ値 $\hat{\theta}$ は、最尤推定量。

$$l(\theta) \text{ を最大にする } \theta \Rightarrow \hat{\theta}$$

これにより決まるモデルは、最尤モデル

$$f(x_1, \dots, x_n | \hat{\theta}) \quad \text{or} \quad p(x_1, \dots, x_n | \hat{\theta})$$

最尤モデルの対数尤度は、最大対数尤度

$$l(\hat{\theta}) = \max_{\theta \in \Theta} l(\theta)$$

B18

最尤法 — 2項分布

$$p(k | \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}; \quad k = 0, 1, \dots, n$$

$$l(\theta) = \log \binom{n}{k} + k \log \theta + (n - k) \log(1 - \theta)$$

これを最大にする θ を求めるために θ で微分すると

$$\frac{dl}{d\theta} = \frac{k}{\theta} - \frac{n-k}{1-\theta}$$

これを0とおいて、

$$\hat{\theta} = \frac{k}{n}$$

B19

最尤法 — 多項分布

c 個の事象の起こった度数を n_1, \dots, n_c , $n_1 + \dots + n_c = n$ とする。

$$l(\theta) = \sum_{i=1}^c n_i \log \theta_i$$

$\sum_{i=1}^c \theta_i = 1$ の条件の下で $l(\theta)$ を最大にする $\hat{\theta}$ は、

$$\hat{\theta}_i = \frac{n_i}{n}$$

として得られる。

$$l(\hat{\theta}) = \sum_{i=1}^c n_i \log n_i - n \log n$$

B20

最尤法 — ポアソン分布

$$\log p(k|\lambda) = -\lambda + k \log \lambda - \log k!$$

i 番目の観測値を k_i とする .

$$l(\lambda) = \sum_{i=1}^n \log p(k_i|\lambda) = \sum_{i=1}^n (-\lambda + k_i \log \lambda - \log k_i!)$$

$$\frac{dl(\lambda)}{d\lambda} = \sum_{i=1}^n \left(-1 + \frac{k_i}{\lambda} \right)$$

これを 0 とおいて

$$\hat{\lambda} = \frac{1}{n} \sum_{i=1}^n k_i$$

B21

最尤法 — 正規分布 $N(\mu, \sigma^2)$

$$l(\mu, \sigma^2) = -\frac{n}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) = 0$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

を解いて ,

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

これを $l(\mu, \sigma^2)$ に代入すると

$$l(\hat{\mu}, \hat{\sigma}^2) = -\frac{n}{2} \log 2\pi\hat{\sigma}^2 - \frac{n}{2}$$

B22

最尤法の性質 (モデルが正しいとき)

最尤推定量の漸近 (ぜんきん) 正規性

$X_1, \dots, X_n \sim f(x|\theta^*)$ なら十分大きな n に対して近似的に

$$\hat{\theta} \sim N(\theta^*, G^{-1}/n)$$

$$G_{ij} = E_X \left[\frac{\partial \log f(X|\theta)}{\partial \theta_i} \frac{\partial \log f(X|\theta)}{\partial \theta_j} \right]_{\theta=\theta^*}$$

$n \rightarrow \infty$ のとき

- $\hat{\theta} \rightarrow \theta^*$ (一致性)
- $E(\hat{\theta}) = \theta^* + O(n^{-1})$ (漸近不偏性)
- 他の推定量では G^{-1}/n より小さな分散には出来ない (漸近有効性)

B23

最尤法の性質

$X_1, \dots, X_n \sim g(x)$ なら十分大きな n に対して近似的に

$$\hat{\theta} \sim N(\theta^*, V/n)$$

$$\theta^* = \arg \min_{\theta \in \Theta} I(g; f(\theta)), \quad V = H^{-1}GH^{-1}$$

$$G_{ij} = E_X \left[\frac{\partial \log f(X|\theta)}{\partial \theta_i} \frac{\partial \log f(X|\theta)}{\partial \theta_j} \right]_{\theta=\theta^*}, \quad H_{ij} = -E_X \left[\frac{\partial^2 \log f(X|\theta)}{\partial \theta_i \partial \theta_j} \right]_{\theta=\theta^*}$$

$f(\cdot|\theta^*)$ はパラメトリック・モデル $f(\cdot|\Theta)$ のなかで最も $g(\cdot)$ に近い分布 .

パラメトリック・モデルが正しい場合 $\Leftrightarrow f(\cdot|\Theta)$ が $g(\cdot)$ を含んでいる

$$g(\cdot) = f(\cdot|\theta^*), \quad H = G = \text{Fisher 情報行列なので } V = G^{-1} = H^{-1}$$

B24

4 AIC

この章のトピック — 赤池情報量規準 **AIC**

- 概要
- 期待平均対数尤度
- **AIC**の導出

B29

モデル選択

複数のパラメトリック・モデルから、良いものを選ぶ。

$$\text{AIC} = -2 \times (\text{最大対数尤度}) + 2 \times (\text{自由パラメタ数})$$

をモデル選択の規準とする

候補となるモデル毎に**AIC**を計算して、その値を最小にするモデルを選択する。
このモデルを、最小**AIC**推定値 (**MAICE**)とよぶ

- **AIC**の第1項: モデルのデータへの当てはまりの良さ
- **AIC**の第2項: モデルの複雑さ

同程度にあてはまりの良いモデルが複数ある場合、そのなかで自由パラメタ数の小さいものを選ぶべきである — 節約の原理

B30

4.1 概要

パラメトリック・モデルの良さの規準

$$\text{期待平均対数尤度} := E\{\text{平均対数尤度}\}$$

期待平均対数尤度を最大対数尤度で推定 — 偏りがある

偏りを補正

$$E\{\text{最大対数尤度} - \text{自由パラメタ数}\} \approx \text{期待平均対数尤度}$$

$$\text{AIC} := -2 \times (\text{最大対数尤度} - \text{自由パラメタ数})$$

B31

定義と仮定

確率変数 $X = (X_1, \dots, X_n)$, データは実現値 $x = (x_1, \dots, x_n)$.

真の分布

$$X_1, \dots, X_n \sim g(\cdot) \quad \text{独立に同分布に従う}$$

$$g(x_1, \dots, x_n) = g(x_1) \cdots g(x_n)$$

モデル **MODEL**(K)

$$X_1, \dots, X_n \sim f(\cdot|\theta) \quad \text{独立に同分布に従う}$$

$$f(x_1, \dots, x_n|\theta) = f(x_1|\theta) \cdots f(x_n|\theta)$$

対数尤度

$$l(\theta) = \sum_{i=1}^n \log f(x_i|\theta); \quad \theta = (\theta_1, \dots, \theta_K)$$

B32

MODEL(K)に制約を入れた, 部分モデル MODEL(k)

$$\Theta_k = \{\theta \in \Theta_K | \theta_{k+1} = \theta_{k+2} = \dots = \theta_K = 0\}$$

MODEL(k)の自由パラメタ数 = k

$$\Theta_1 \subset \Theta_2 \subset \dots \subset \Theta_K$$

データ x が与えられたときの MODEL(k) のパラメタの最尤推定量を $\hat{\theta}_k$ とする .

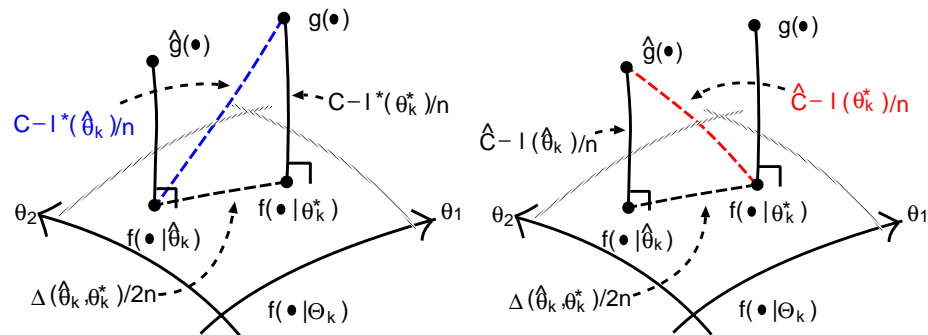
最大対数尤度

$$l(\hat{\theta}_k) = \max_{\theta \in \Theta_k} l(\theta)$$

情報量規準

$$AIC(k) = -2 \times (l(\hat{\theta}_k) - k)$$

B33



$$E_X \{ I(g; f(\hat{\theta}_k)) \} = C - l_n^*(k)/n$$

$$I(g; f(\hat{\theta}_k)) - I(g; f(\theta_k^*)) \approx \Delta(\hat{\theta}_k, \theta_k^*)/2n$$

$$I(\hat{g}; f(\theta_k^*)) - I(\hat{g}; f(\hat{\theta}_k)) \approx \Delta(\hat{\theta}_k, \theta_k^*)/2n$$

$$E_X \{ l(\theta_k^*) \} = l^*(\theta_k^*)$$

B35

4.2 期待平均対数尤度

分布 $f(\cdot|\theta)$ の平均対数尤度は ,

$$E_Z \{ \log f(Z|\theta) \} = \int g(z) \log f(z|\theta) dz; \quad Z \sim g(\cdot)$$

これを n 倍して

$$l^*(\theta) = n E_Z \{ \log f(Z|\theta) \}$$

$l^*(\theta)$ が大きい程 $f(\cdot|\theta)$ の $g(\cdot)$ に対する近似が良いことになる .

MODEL(k) の最尤モデルの良さは $f(\cdot|\hat{\theta}_k)$ の平均対数尤度 $l^*(\hat{\theta}_k)$ で評価

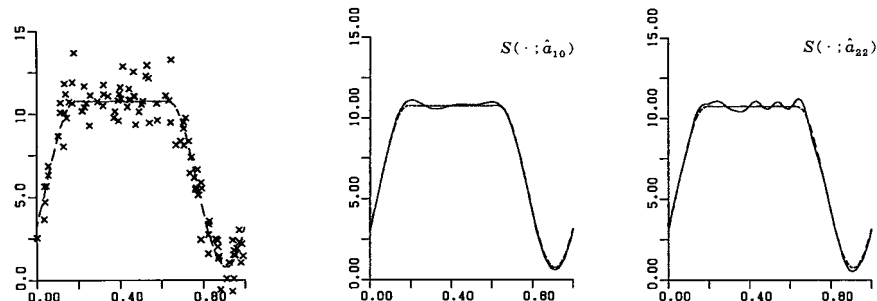
$l^*(\hat{\theta}_k)$ のデータに関する期待値を取って「期待平均対数尤度」

$$l_n^*(k) = E_X \{ l^*(\hat{\theta}_k) \} = \int l^*(\hat{\theta}_k) g(x) dx$$

でパラメトリック・モデル MODEL(k) を評価する

B34

【数値例】 条件つき分布モデル (フーリエ級数による回帰モデル)



K = 22 のモデルから
データを生成

k = 10 と k = 22 のモデルの当てはめ

B36

MODEL(K):

$$Y \sim N(S(x; \mathbf{a}), \sigma^2), \quad X \sim [0, 1] \text{ の一様分布}$$

$$S(x; \mathbf{a}) = a_0 + \sum_{m=1}^M (a_{2m-1} \sin 2m\pi x + a_{2m} \cos 2m\pi x)$$

$$f(x, y|\theta) = \begin{cases} \exp(-(y - S(x; \mathbf{a}))^2/2\sigma^2)/\sqrt{2\pi\sigma^2} & 0 \leq x \leq 1 \\ 0 & x < 0, x > 1 \end{cases}$$

$$\theta \in \Theta = \{(\theta_1, \dots, \theta_K) = (\sigma^2, a_0, \dots, a_{2M}) | \sigma^2 > 0\}$$

$$K = 2M + 2$$

$$\Theta_1 \subset \Theta_2 \subset \dots \subset \Theta_K$$

B37

表4.1の最後の2行

$$l(\hat{\theta}_2) \leq l(\hat{\theta}_4) \leq \dots \leq l(\hat{\theta}_K)$$

最大対数尤度 $l(\hat{\theta}_k)$ は k とともに大きくなる

$$l_{500}^*(2) \leq l_{500}^*(4) \leq \dots \leq l_{500}^*(10) \geq l_{500}^*(12) \geq \dots \geq l_{500}^*(22)$$

期待平均対数尤度 $l_n^*(k)$ は $k = 10$ で最大値をとる

真のモデルは $k = 22$ は必ずしも期待平均対数尤度を大きくしない

B39

表 4.1 最尤推定値, 最大対数尤度および期待平均対数尤度

θ	$\hat{\theta}_2$	$\hat{\theta}_4$	$\hat{\theta}_6$	$\hat{\theta}_8$	$\hat{\theta}_{10}$	$\hat{\theta}_{12}$	$\hat{\theta}_{14}$	$\hat{\theta}_{16}$	$\hat{\theta}_{18}$	$\hat{\theta}_{20}$	$\hat{\theta}_{22}$	θ^*
θ_1	14.220	4.145	1.121	0.984	0.974	0.970	0.968	0.966	0.956	0.955	0.949	1.000
θ_2	8.012	7.983	8.043	8.010	8.010	8.012	8.013	8.013	8.014	8.013	8.011	8.000
θ_3	2.470	2.491	2.488	2.484	2.484	2.479	2.479	2.479	2.470	2.468	2.465	2.415
θ_4	-3.684	-3.869	-3.854	-3.850	-3.851	-3.853	-3.852	-3.847	-3.847	-3.847	-3.847	-3.806
θ_5		2.252	2.238	2.242	2.248	2.248	2.247	2.247	2.247	2.253	2.249	2.119
θ_6		-0.984	-0.987	-0.987	-0.995	-0.994	-0.997	-1.005	-1.006	-1.006	-1.009	-0.997
θ_7			0.523	0.516	0.515	0.518	0.518	0.512	0.515	0.515	0.515	0.545
θ_8			-0.022	-0.013	-0.012	-0.010	-0.007	0.002	0.004	0.007	0.007	0.069
θ_9			-0.129	-0.128	-0.130	-0.126	-0.133	-0.136	-0.139	-0.139	-0.139	-0.094
θ_{10}			-0.074	-0.074	-0.074	-0.073	-0.072	-0.071	-0.073	-0.073	-0.073	-0.078
θ_{11}				0.070	0.076	0.079	0.079	0.081	0.089	0.089	0.089	-0.021
θ_{12}				0.041	0.041	0.038	0.024	0.024	0.021	0.021	0.021	-0.065
θ_{13}				-0.058	-0.058	-0.061	-0.061	-0.061	-0.060	-0.060	-0.060	-0.011
θ_{14}				0.032	0.036	0.042	0.047	0.044	0.044	0.044	0.044	0.042
θ_{15}				-0.065	-0.067	-0.069	-0.075	-0.075	-0.075	-0.075	-0.075	-0.011
θ_{16}				0.003	-0.010	-0.010	-0.011	-0.011	-0.011	-0.011	-0.011	0.010
θ_{17}					0.145	0.146	0.142	0.142	0.142	0.142	0.142	0.020
θ_{18}					-0.024	-0.024	-0.016	-0.016	-0.016	-0.016	-0.016	-0.004
θ_{19}					-0.034	-0.035	-0.035	-0.035	-0.035	-0.035	-0.035	0.002
θ_{20}					0.041	0.039	0.039	0.039	0.039	0.039	0.039	0.001
θ_{21}						0.050	0.050	0.050	0.050	0.050	0.050	-0.006
θ_{22}						-0.089	-0.089	-0.089	-0.089	-0.089	-0.089	-0.009
$-l(\hat{\theta}_k)$	1373.14	1064.97	738.13	705.46	702.93	701.91	701.41	700.90	698.35	697.84	696.30	
$-l_{500}^*(k)$	1371.11	1051.87	750.07	716.43	715.73	716.28	717.20	718.29	719.39	720.57	721.80	

B38

数値例における最大対数尤度と期待平均対数尤度の関係

表 4.2 1000回の実験で求めた最大対数尤度の平均値および $l_{500}^*(k)$ との差

	2	4	6	8	10	12
$l(\hat{\theta}_k)$ の平均	-1368.85	-1047.52	-744.53	-708.55	-705.61	-703.94
$l(\hat{\theta}_k) - l_{500}^*(k)$ の平均	2.26	4.35	5.54	7.88	10.12	12.34
	14	16	18	20	22	
$l(\hat{\theta}_k)$ の平均	-702.68	-701.65	-700.61	-699.55	-698.49	
$l(\hat{\theta}_k) - l_{500}^*(k)$ の平均	14.52	16.64	18.75	21.02	23.31	

$$E \{l(\hat{\theta}_k) - l_n^*(k)\} \approx k$$

B40

4.3 AICの導出

[仮定] 以下の議論では真の分布を

$$g(\cdot) = f(\cdot|\theta^*), \quad \theta^* = (\theta_1^*, \dots, \theta_K^*)$$

とおく. $\theta^* \in \Theta_k$ となる最小の k を真の自由パラメータ数 k^* とする.

$$\theta^* = (\theta_1^*, \dots, \theta_k^*, 0, \dots, 0)$$

[情報量規準]

$$\text{AIC}(k) = -2 \times (l(\hat{\theta}_k) - k)$$

[導出する式]

$$E \{ l(\hat{\theta}_k) - k \} \approx l_n^*(k)$$

B41

第1段階

平均対数尤度 (n 倍)

$$l^*(\theta) = nE_Z \{ \log f(Z|\theta) \}$$

を θ^* の周辺でテーラ展開して近似式

$$\begin{aligned} l^*(\theta) &\approx l^*(\theta^*) + n(\theta - \theta^*)E_Z \left[\frac{\partial \log f(Z|\theta)}{\partial \theta'} \right]_{\theta^*} \\ &\quad + \frac{1}{2}n(\theta - \theta^*)E_Z \left[\frac{\partial^2 \log f(Z|\theta)}{\partial \theta' \partial \theta} \right]_{\theta^*} (\theta - \theta^*)' \\ &= l^*(\theta^*) - \frac{1}{2}\sqrt{n}(\theta - \theta^*)H\sqrt{n}(\theta - \theta^*)' \\ &= l^*(\theta^*) - \frac{1}{2}\Delta(\theta, \theta^*) \end{aligned}$$

ただし

$$\Delta(\theta, \theta^*) = \sqrt{n}(\theta - \theta^*)H\sqrt{n}(\theta - \theta^*)'$$

B42

MODEL(K)での最尤推定 $\hat{\theta}_K$ を $l^*(\theta)$ に代入すると

$$l^*(\hat{\theta}_K) \approx l^*(\theta^*) - \frac{1}{2}\Delta(\hat{\theta}_K, \theta^*)$$

ただし

$$\hat{\theta}_K \sim N(\theta^*, V/n)$$

ここで $H = G$, $V = H^{-1}GH^{-1} = H^{-1}$ を利用すると,

$$\Delta(\hat{\theta}_K, \theta^*) \sim \chi_K^2$$

$E\{\chi_K^2\} = K$ なので

$$l_n^*(K) = E_X \{ l^*(\hat{\theta}_K) \} \approx l^*(\theta^*) - \frac{K}{2}$$

B43

MODEL(k)においては $H_k \approx G_k$ がいえて, やはり

$$l_n^*(k) = E_X \{ l^*(\hat{\theta}_k) \} \approx l^*(\theta_k^*) - \frac{k}{2}$$

もうすこし厳密にやると,

$$\begin{aligned} l_n^*(k) &\approx l^*(\theta_k^*) - \frac{1}{2} \text{tr} \left(H_k E_X \{ n(\hat{\theta}_k - \theta_k^*)'(\hat{\theta}_k - \theta_k^*) \} \right) \\ &\approx l^*(\theta_k^*) - \frac{1}{2} \text{tr} \left(H_k H_k^{-1} G_k H_k^{-1} \right) \\ &= l^*(\theta_k^*) - \frac{1}{2} \text{tr} \left(G_k H_k^{-1} \right) \\ &\quad \text{tr} \left(G_k H_k^{-1} \right) \approx k \end{aligned}$$

$$G_k = E_X \left[\frac{\partial \log f(X|\theta)}{\partial \theta'} \frac{\partial \log f(X|\theta)}{\partial \theta} \right]_{\theta=\theta_k^*} \quad H_k = -E_X \left[\frac{\partial^2 \log f(X|\theta)}{\partial \theta' \partial \theta} \right]_{\theta=\theta_k^*}$$

ただしここでの計算では θ は最初の k 個の要素に制限してある

B44

第2段階

対数尤度

$$l(\theta) = \sum_{i=1}^n \log f(x_i|\theta)$$

を $\hat{\theta}_K$ の周辺でテーラ展開すると,

$$l(\theta) \approx l(\hat{\theta}_K) + (\theta - \hat{\theta}_K) \left[\frac{\partial l}{\partial \theta'} \right]_{\hat{\theta}_K} + \frac{1}{2} (\theta - \hat{\theta}_K) \left[\frac{\partial^2 l}{\partial \theta' \partial \theta} \right]_{\hat{\theta}_K} (\theta - \hat{\theta}_K)'$$

ここで

$$\frac{1}{n} \left[\frac{\partial^2 l}{\partial \theta' \partial \theta} \right]_{\hat{\theta}_K} = \frac{1}{n} \sum_{i=1}^n \left[\frac{\partial^2 \log f(x_i|\theta)}{\partial \theta' \partial \theta} \right]_{\hat{\theta}_K} \rightarrow E_Z \left[\frac{\partial^2 \log f(Z|\theta)}{\partial \theta' \partial \theta} \right]_{\theta^*} = -H$$

を代入すると

$$l(\theta) \approx l(\hat{\theta}_K) - \frac{1}{2} \Delta(\theta, \hat{\theta}_K)$$

B45

MODEL(K)での最適なパラメタ θ^* を $l(\theta)$ に代入すると

$$l(\theta^*) \approx l(\hat{\theta}_K) - \frac{1}{2} \Delta(\theta^*, \hat{\theta}_K)$$

これに次の二つの式を代入して

$$E_X \{l(\theta^*)\} = E_X \left\{ \sum_{i=1}^n \log f(X_i|\theta^*) \right\} = n E_Z \{ \log f(Z|\theta^*) \} = l^*(\theta^*)$$

$$E_X \{ \Delta(\theta^*, \hat{\theta}_K) \} \approx K$$

MODEL(K)では

$$l^*(\theta^*) \approx E_X \{ l(\hat{\theta}_K) \} - \frac{K}{2}$$

同様に MODEL(k)では

$$l^*(\theta_k^*) \approx E_X \{ l(\hat{\theta}_k) \} - \frac{k}{2}$$

B46

第1段階と第2段階をまとめる

$$l_n^*(k) \approx l^*(\theta_k^*) - \frac{k}{2}$$

$$l^*(\theta_k^*) \approx E_X \{ l(\hat{\theta}_k) \} - \frac{k}{2}$$

MODEL(k)

$$l_n^*(k) \approx E_X \{ l(\hat{\theta}_k) \} - k$$

もうすこし厳密にやると,

$$l_n^*(k) \approx E_X \{ l(\hat{\theta}_k) \} - \text{tr}(G_k H_k^{-1})$$

赤池の情報量規準

$$\text{AIC}(k) = -2 \times (l(\hat{\theta}_k) - k)$$

竹内の情報量規準

$$\text{TIC}(k) = -2 \times (l(\hat{\theta}_k) - \text{tr}(G_k H_k^{-1}))$$

B47

AICのバラツキ

表 4.6 AIC 最小となる $k (=k_A)$ の分布 ($n = 500$)

k_A	2	4	6	8	10	12	14	16	18	20	22	計
度数	0	0	0	326	342	173	91	27	17	14	10	1000

表 4.7 $n = 2000$ のときの k_A の分布

k_A	2	4	6	8	10	12	14	16	18	20	22	計
度数	0	0	0	6	210	359	264	68	58	21	14	1000

なお, $l_{500}^*(k)$ 最小は $k = 10$, $l_{2000}^*(k)$ 最小は $k = 12$ のとき

B48

AIC 利用上の注意

- データ数に対して自由パラメタ数が大きすぎる場合は **AIC** は使えない
(導出の際に, k を固定して $n \rightarrow \infty$ という漸近論を用いている)
- **AIC** の値そのものよりも, **AIC** の値の差に意味がある
- モデル比較の際は, **AIC** の差のバラツキを考慮する
とくに「ノンネスト」の場合は, バラツキが非常に大きい!
- モデル比較で **AIC** の値がほぼ等しくても推定された分布の形が大幅に違うこともある
- モデルの次数を上げるほど **AIC** が下がってなかなか最小値を取らないことがある
- **AIC** は真の自由パラメタ数を推定するための規準ではない!